

The QCDOC Project Overview and Status

DOE LGT Review

May 24-25, 2005

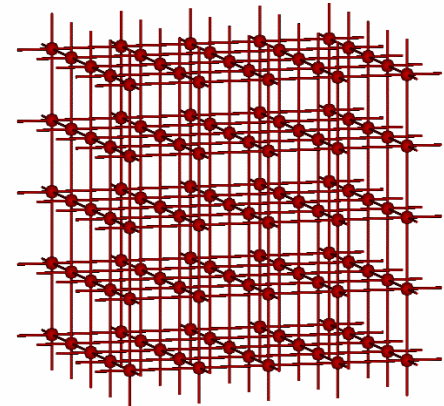
Norman H. Christ

Outline

- Project goals
- QCDOC collaboration
- Architecture
- Software
 - Operating system
 - Run-time environment
 - Programming environment
- Construction and packaging
- Construction status
- Final bring-up issues
- Application performance
- Future plan

Project Goals

- Lattice QCD provides the only first-principles window into non-perturbative phenomena of QCD.
- All significant errors are controlled and can be reduced with faster computers or better algorithms.
- Simple formulation enables targeted computer architecture.
- Regular space-time description: easily mounted on a parallel computer.



8192-node, 0.4 Tflops QCDS machine

Project Goals (con't)

- Massively parallel machine capable of *strong scaling*: use many nodes on a small problem.
 - Large inter-node bandwidth.
 - Small communications latency.
- \$1/sustained Mflops cost/performance.
- Low power, easily maintained modular design.

QCDOC Collaboration

(people)

- Columbia (DOE)
 - Norman Christ
 - Saul Cohen*
 - Calin Cristian*
 - Zhihua Dong
 - Changhoan Kim*
 - Ludmila Levkova*
 - Sam Li*
 - Xiaodong Liao*
 - Guofeng Liu*
 - Meifeng Lin*
 - Robert Mawhinney
 - Azusa Yamaguchi
- BNL (SciDAC)
 - Robert Bennett
 - Chulwoo Jung
 - Konstantin Petrov
 - David Stampf
- UKQCD (PPARC)
 - Peter Boyle
 - Mike Clark
 - Balint Joo
- RBRC (RIKEN)
 - Shigemi Ohta
 - Tilo Wettig
- IBM
 - Dong Chen
 - Alan Gara
 - Design groups:
 - Yorktown Heights, NY
 - Rochester, MN
 - Raleigh, NC

*CU graduate student

QCDOC Collaboration

(money)

Institution/funding source	Design and proto-typing	Large installations
Columbia/DOE	\$500K	\$1M (UKQCD)
RBRC/RIKEN	\$400K	\$5M
UKQCD/PPARC	\$1M	\$5.2M
BNL/DOE	-	\$5.1M

Personnel and site prep costs are not included.

QCDOC Architecture

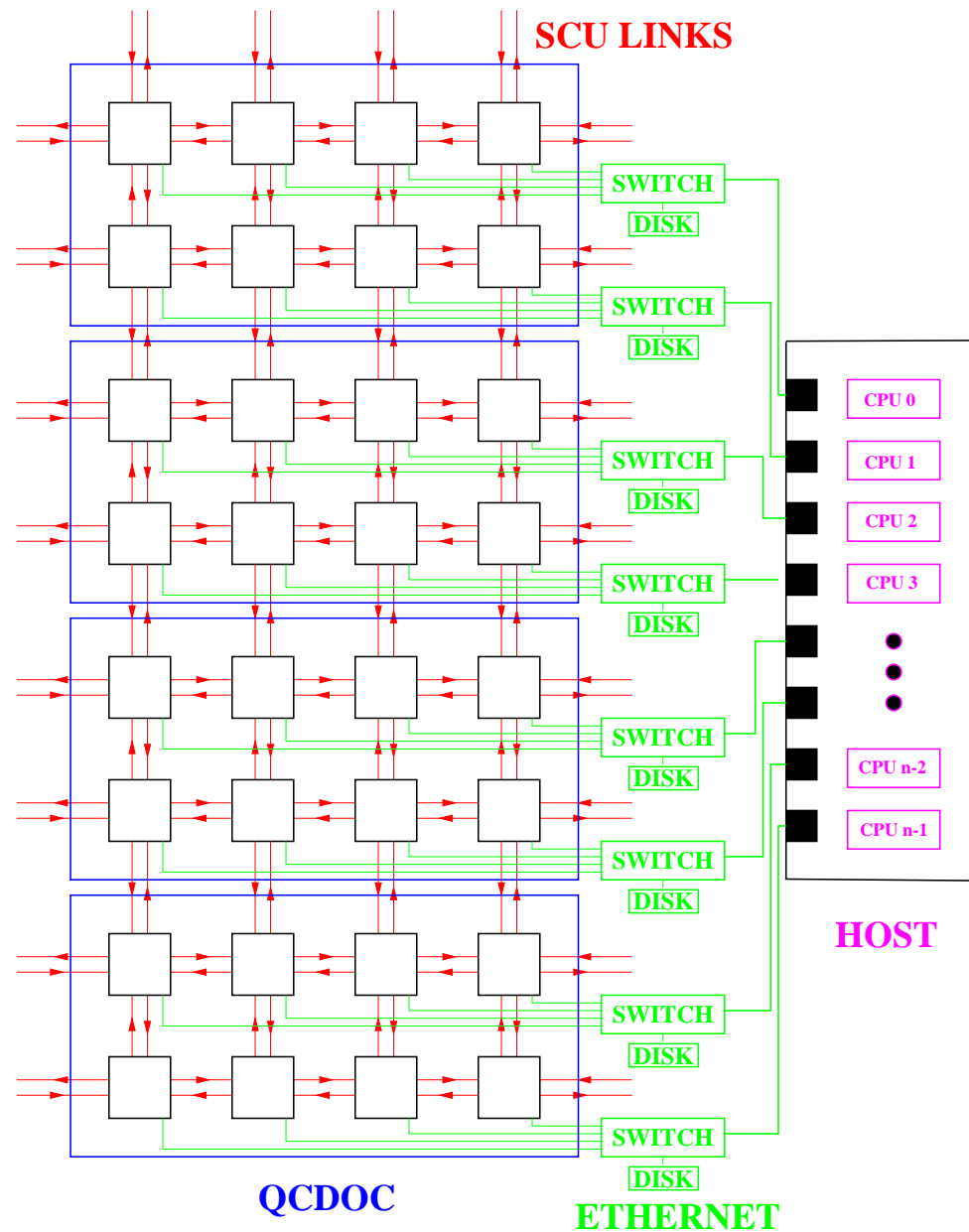
- IBM-fabricated, single-chip node.
[50 million transistors, 5 Watt, 1.3cm x 1.3cm]
- Processor:
 - PowerPC 32-bit RISC.
 - 64-bit, 1 Gflops floating point unit.
- Memory/node: 4 Mbyte (on-chip) & 2 Gbyte DIMM.
- Communications network:
 - 6-dim, supporting lower dimensional partitions.
 - Global sum/broadcast functionality.
 - Multiple DMA engines/minimal processor overhead.
- Ethernet connection to each node: booting, I/O, host control.
- ~7-8 Watt/node, 15 in³ per node.

Software Environment

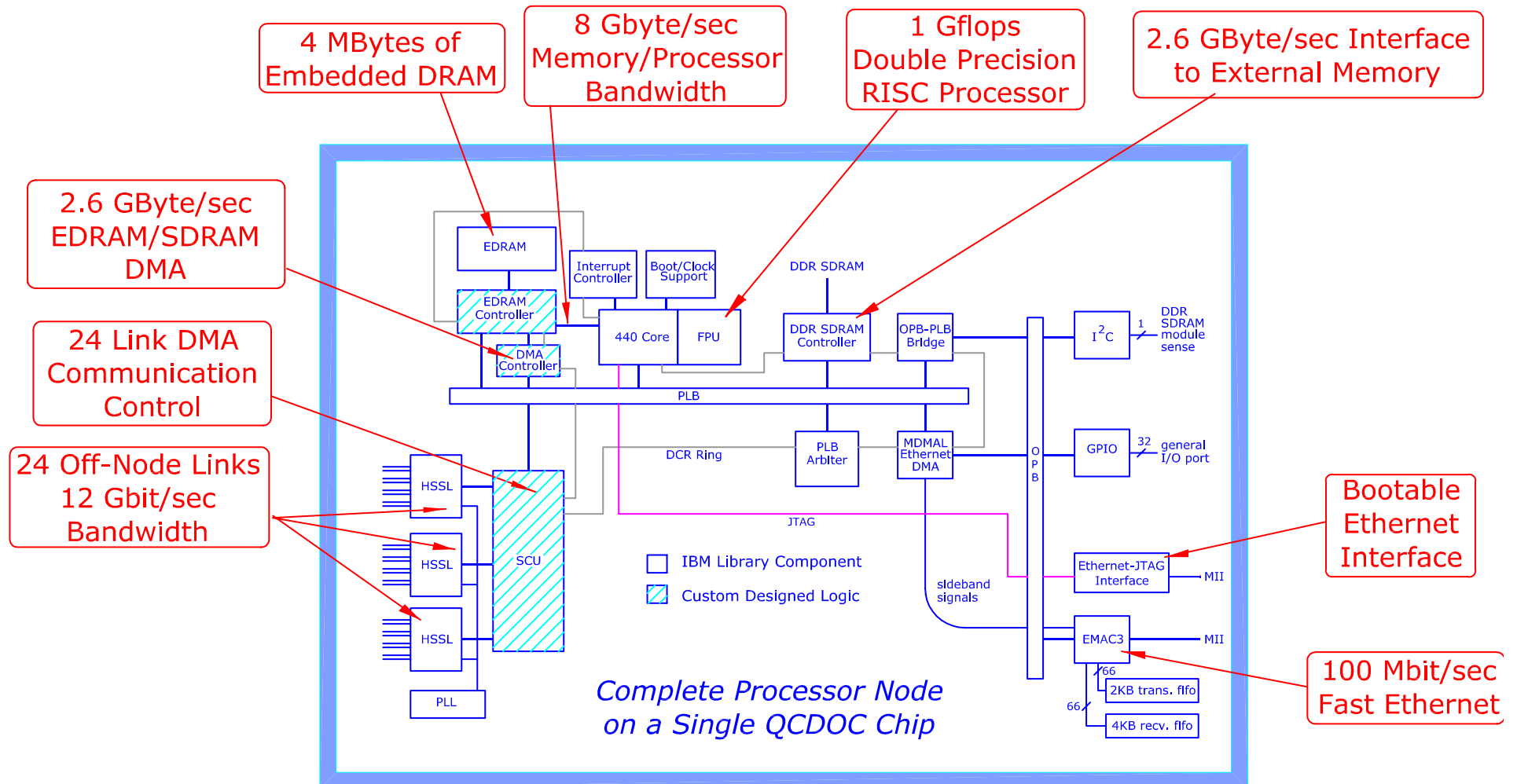
- Lean kernel on each node
 - Protected kernel mode and address space.
 - RPC support for host access.
 - NFS access to NAS disks (/pfs).
 - Normal Unix services including stdout and stderr.
- Threaded host kernel
 - Efficient performance on 8-processor SMP host.
 - User shell (qsh) with extended commands.
 - Host file system (/host).
 - Simple remapping of 6-D machine to (6-n)-D torus.
- Programming environment
 - POSIX compatible, open-source libc.
 - gcc and xlc compilers
- SciDAC standards
 - Level-1, QMP protocol
 - Level-2 parallelized linear algebra, QDP & QDP++.
 - Efficient level-3 inverters
 - Wilson/clover
 - Domain wall fermions
 - ASQTAD
 - p4 (underway)

Network Architecture

- Red boxes are nodes.
- Blue boxes mother boards.
- Red lines are communications links.
- Green lines are Ethernet connections.
- Green boxes are Ethernet switches.
- Pink boxes are host CPU processors.

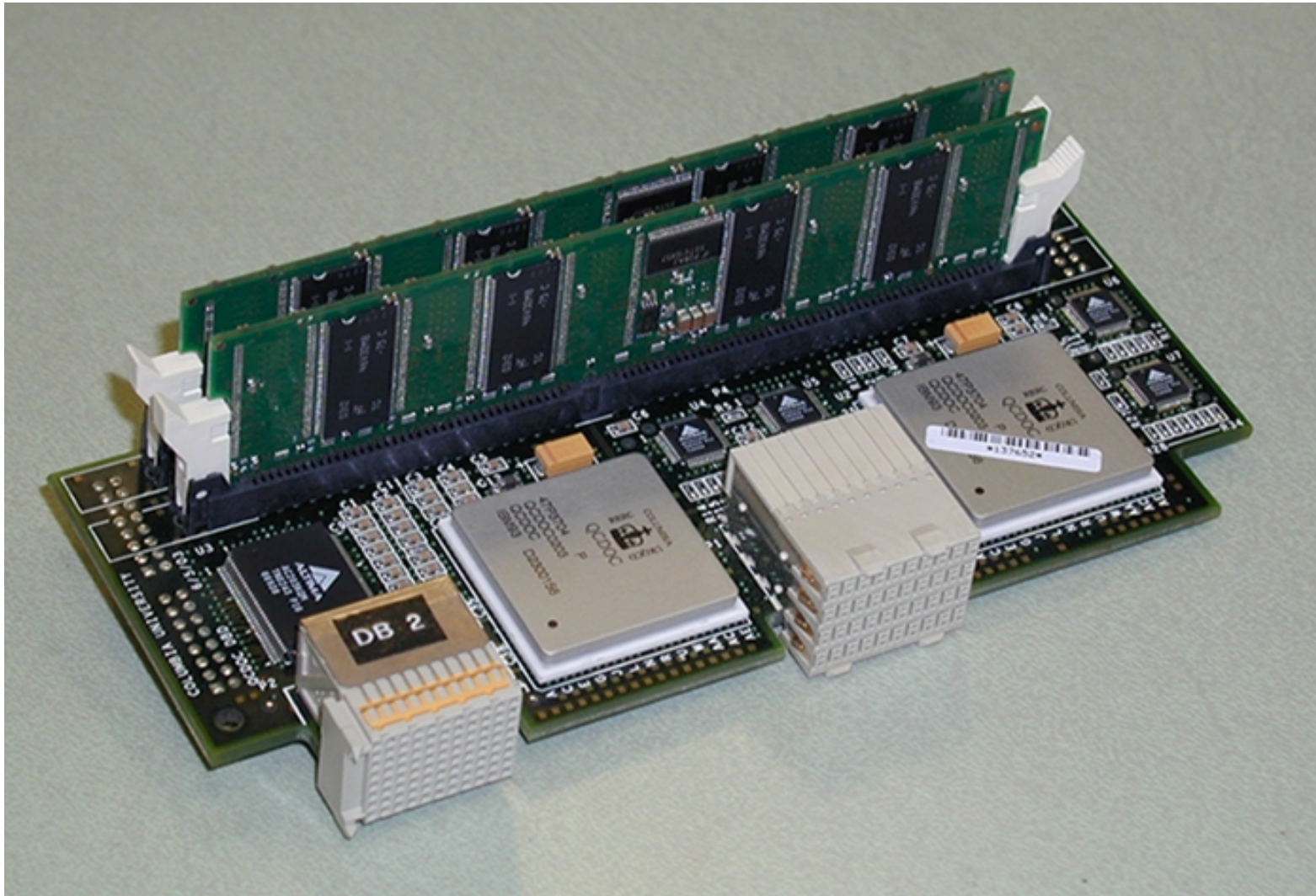


QCDOC Chip



50 million transistors, 0.18 micron, 1.3 x 1.3 cm die, 5 Watt

Daughter board (2 nodes)



Mother board (64 nodes)



512-Node Machine



UKQCD Machine (12,288 nodes/10 Tflops)



Brookhaven Installation



RBRC (right) and DOE (left) 12K-node QCDOC machines

Project Status

- UKQCD – 13,312 nodes --\$5.2M – 3-5 Tflops sustained.
 - Installed in Edinburgh 12/04.
 - Running production at 400 MHz/100% reprod.
- RBRC – 12,288 nodes -- \$5M – 3-5 Tflops sustained.
 - Installed at BNL 2/05.
 - 1/3 in production/100% reprod.
 - 1/3 performing physics tests.
 - 1/3 speed sorting 420.
- DOE – 12,288 nodes -- \$5.1M – 3-5 Tflops sustained
 - Installed at BNL 4/05.
 - 1/2 performing physics tests.
 - 1/2 being debugged.
- Price/performance of **~\$1/Mflops.**

Final Bring-up issues

- FPU errors
 - Lowest two bits infrequently incorrect (not seen at 400MHz).
 - Remove slow nodes at 432MHz and run at 400MHz.
- Serial communication errors.
 - Induced by Ethernet activity.
 - 0.25/month at 400 MHz/1K nodes.
 - Further reduced by PLL tuning.
 - Protected by hardware checksums with no performance loss.
- Parallel disk system
 - 24 Tbyte RAID servers.
 - 512-nodes achieve 12 Mbytes/sec.
 - Installed 05/05?
- Larger machine partitions
 - Three 4096-node partitions assembled.
 - Expect to run as 4096 + 8192 node machines.
- Spares
 - 1% non-functioning daughter boards
 - 1.5% non-functioning mother boards
 - ~18 mother boards for small jobs/code development.

Application Performance

(double precision)

1024-node machine:

Fermion action	Local volume	Dirac performance	CG performance
Wilson	2^4	44%	32%
Wilson	4^4	44%	38%
Clover	4^4	54%	47.5%
DWF	4^5	47%	42%
ASQTAD	4^4	42%	40%

4096-node machine (UKQCD):

DWF/ $24^3 \times 64$ /RHMC (Local vol: $6 \times 6 \times 6 \times 2 \times 8$)	CG:	1.1 Tflops (34%)
	Complete code:	0.95 Tflops (29%)

QCDOC Summary

- Present DOE QCDOC machine use:
 - Alpha users developing code on 1 mother board machines.
 - MILC (staggered 2+1 flavor) using 1K-node machine.
 - JLAB (DWF 2+1 flavor) using 1K-node machine.
 - RBC/BNL (QCD thermo) using 2, 1K-node machines.
 - RBC (DWF) using 1K-node machine.
- 4K node machine being debugged for MILC use.
- Most of machine in production by early June?

The Future: QCDOC++

- Reduced feature size and increased integration permits many nodes per chip (multi-core trend).
- QCDOC → QCDOC++
 - Clock speed (GHz): 0.4 → 1?
 - Integration (nodes/chip): 1 → 64?
 - Performance (Gflops): 0.4 → 2?
 - Inter-node comms (Gbyte/sec): 1 → 10?
 - On chip memory (Mbytes/chip): 4 → 32?
- **Target:** \$0.01/Mflops price performance (1/100 x QCDOC).
- 100x speed-up permits increased 2x problem size per chip.
- Design starts 2006 (with off-project support).
- Target is ambitious with risk. May provide a candidate production machine in 2009.